

Topic 6: - Topic Analysis

Clarissa Boyajian

2022-05-10

```
## -- read in, clean, and wrangle data -- ##
comments_df<-read_csv(here("data/comments_df.csv"))

epa_corp <- corpus(x = comments_df, text_field = "text")
epa_corp.stats <- summary(epa_corp)
tokens <- tokens(epa_corp, remove_punct = TRUE, remove_numbers = TRUE)

# project-specific stop words
add_stops <- c(stopwords("en"),"environmental", "justice", "ej", "epa", "public", "comment")
tokens_1 <- tokens_select(tokens, pattern = add_stops, selection = "remove")

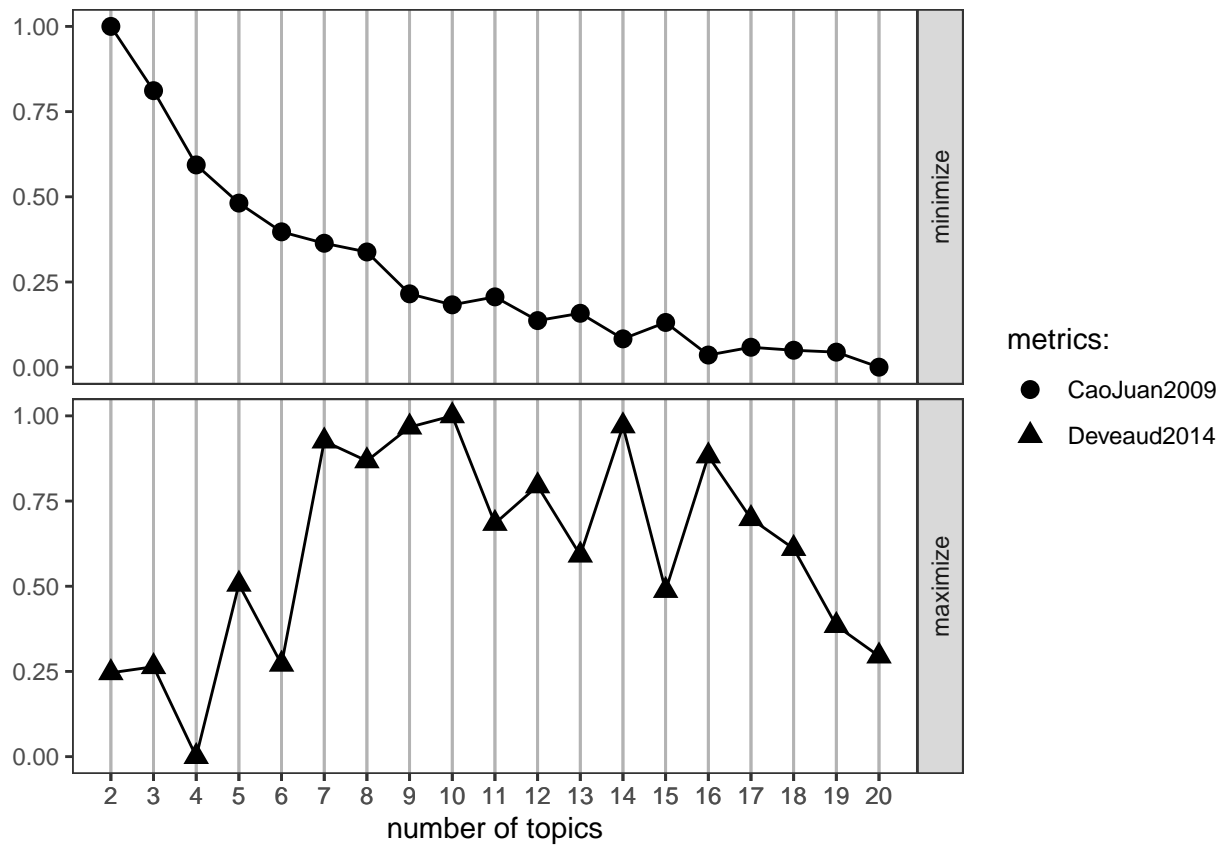
doc_feature_matrix_common <- dfm(tokens_1, tolower = TRUE)
doc_feature_matrix <- dfm_wordstem(doc_feature_matrix_common)
doc_feature_matrix <- dfm_trim(doc_feature_matrix, min_docfreq = 2)

# remove rows (docs) with all zeros
select_idx <- slam::row_sums(doc_feature_matrix) > 0
doc_feature_matrix <- doc_feature_matrix[select_idx, ]
```

Calculate value of k that is most likely

```
# calculate what initial value of k is most likely
result <- FindTopicsNumber(
  doc_feature_matrix,
  topics = seq(from = 2, to = 20, by = 1),
  metrics = c("CaoJuan2009", "Deveaud2014"),
  method = "Gibbs",
  control = list(seed = 77),
  verbose = TRUE
)
```

```
FindTopicsNumber_plot(result)
```



Model 1 ($k = 10$)

```
k <- 10

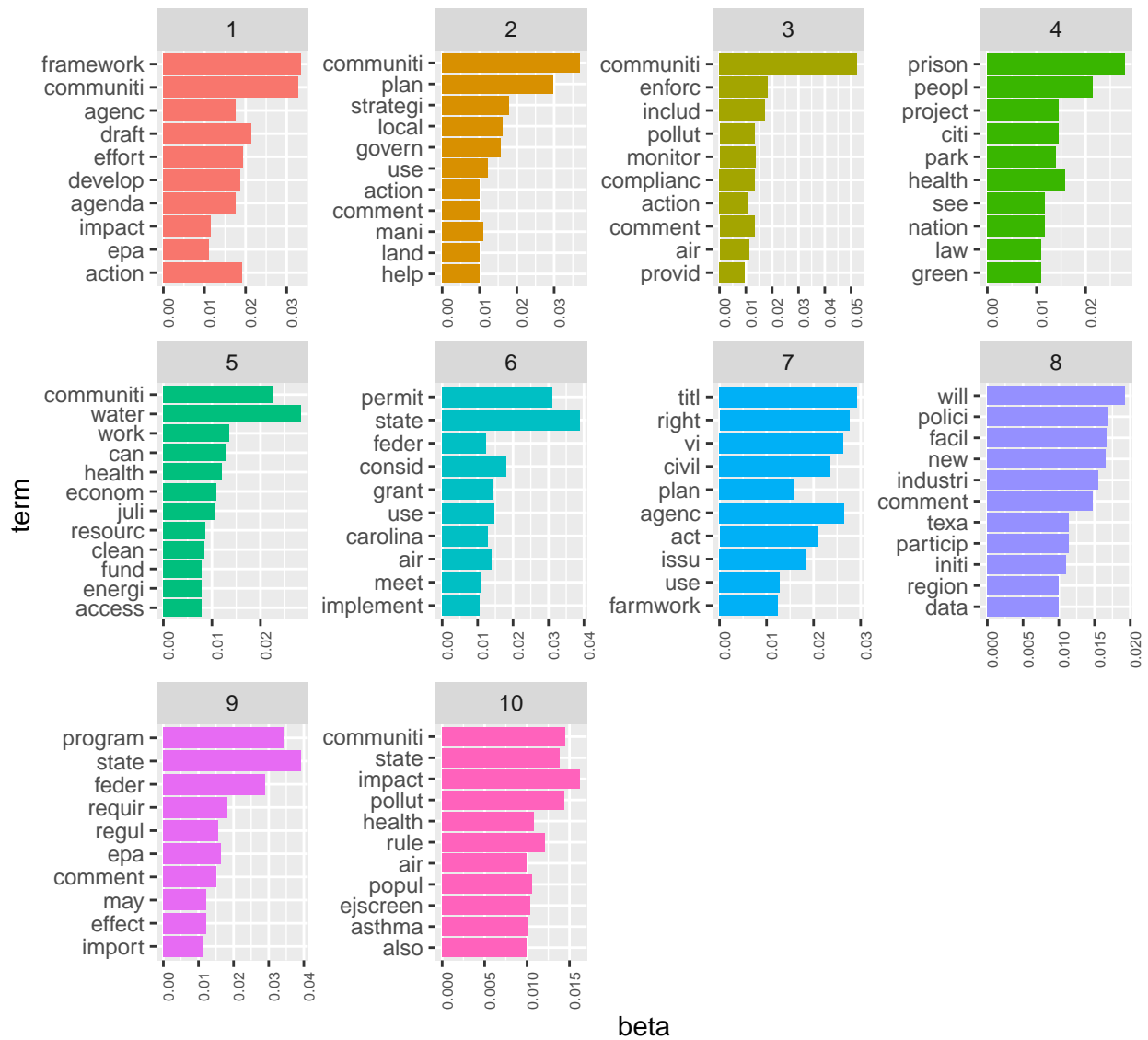
topicModel_k10 <- LDA(doc_feature_matric, k,
  method = "Gibbs",
  control = list(iter = 500, verbose = 25))

tmResult_k10 <- posterior(topicModel_k10)
theta_k10 <- tmResult_k10$topics
beta_k10 <- tmResult_k10$terms # probability of each term in each topic
vocab_k10 <- (colnames(beta_k10))

comment_topic_k10 <- tidy(topicModel_k10, matrix = "beta")

top_terms_k10 <- comment_topic_k10 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)
```

```
top_terms_k10 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



Model 2 ($k = 14$)

```
k <- 14

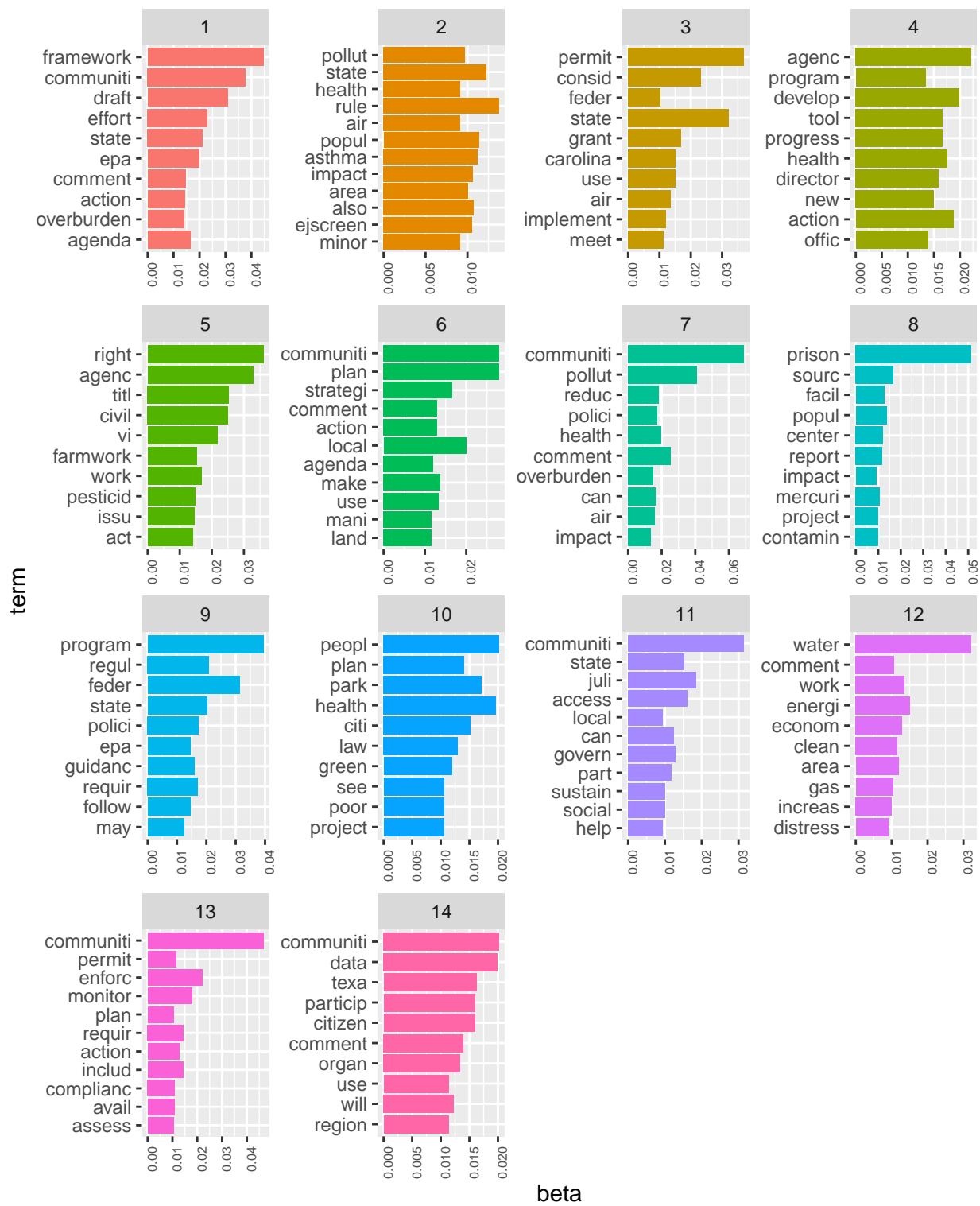
topicModel_k14 <- LDA(doc_feature_matric, k,
                      method = "Gibbs",
                      control = list(iter = 500, verbose = 25))

tmResult_k14 <- posterior(topicModel_k14)
theta_k14 <- tmResult_k14$topics
beta_k14 <- tmResult_k14$terms # probability of each term in each topic
vocab_k14 <- (colnames(beta_k14))

comment_topic_k14 <- tidy(topicModel_k14, matrix = "beta")

top_terms_k14 <- comment_topic_k14 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms_k14 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



Model 3 (k = 16)

```
k <- 16

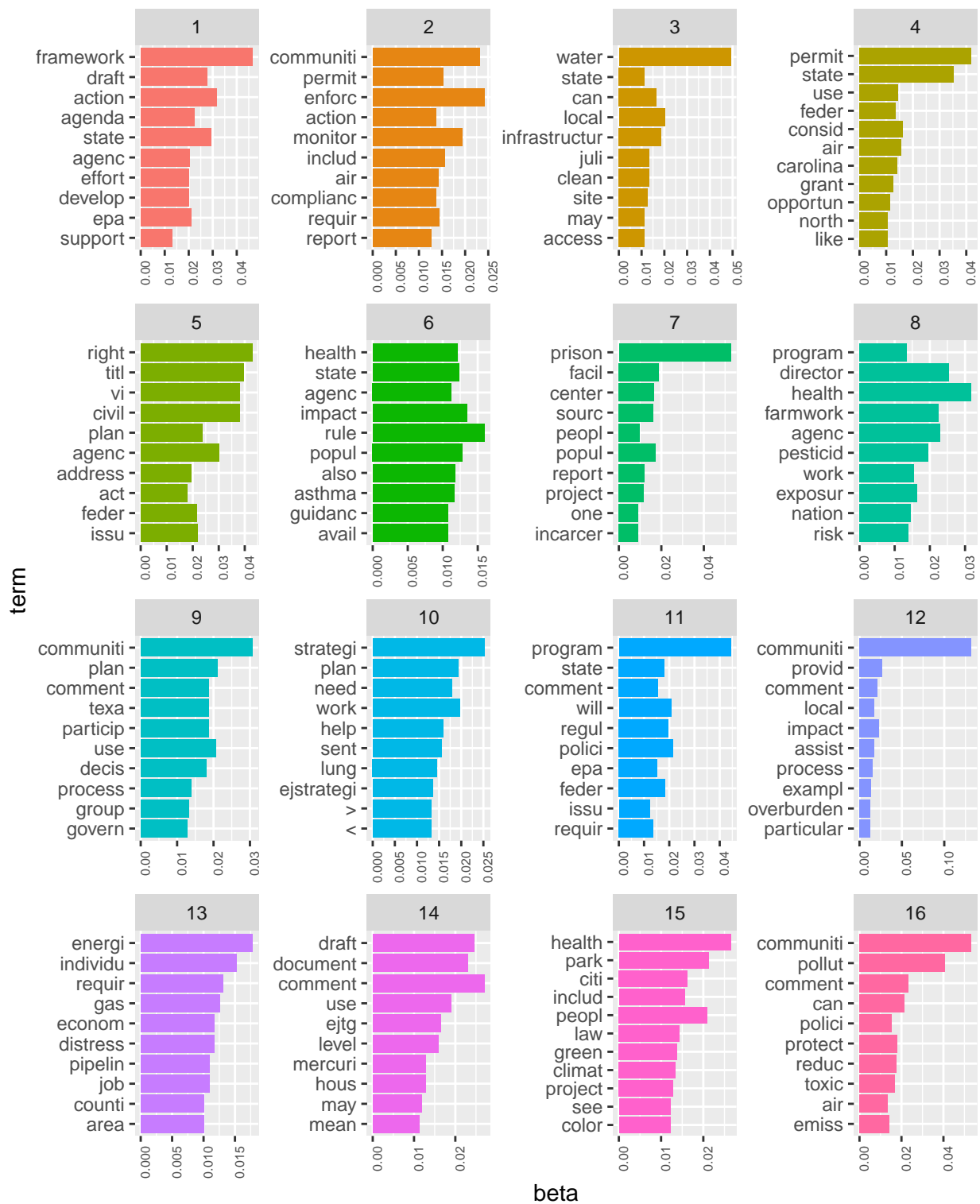
topicModel_k16 <- LDA(doc_feature_matric, k,
                      method = "Gibbs",
                      control = list(iter = 500, verbose = 25))

tmResult_k16 <- posterior(topicModel_k16)
theta_k16 <- tmResult_k16$topics
beta_k16 <- tmResult_k16$terms # probability of each term in each topic
vocab_k16 <- (colnames(beta_k16))

comment_topic_k16 <- tidy(topicModel_k16, matrix = "beta")

top_terms_k16 <- comment_topic_k16 %>%
  group_by(topic) %>%
  top_n(10, beta) %>%
  ungroup() %>%
  arrange(topic, -beta)

top_terms_k16 %>%
  mutate(term = reorder(term, beta)) %>%
  ggplot(aes(term, beta, fill = factor(topic))) +
  geom_col(show.legend = FALSE) +
  facet_wrap(~ topic, scales = "free") +
  coord_flip() +
  theme(axis.text.x = element_text(angle = 90, hjust = 1, size = 6))
```



Evaluate best k value

```
docIds <- c(1, 2, 3, 4, 5, 6)
N <- length(docIds)

## -- k = 10 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k10 <- terms(topicModel_k10, 5)
topicNames_k10 <- apply(top5termsPerTopic_k10, 2, paste, collapse = " ")

# get topic proportions from example documents
topicProportions_K10 <- theta_k10[docIds,]
colnames(topicProportions_K10) <- topicNames_k10

vizDataFrame_k10 <- melt(cbind(data.frame(topicProportions_K10),
                                document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k10 <- ggplot(data = vizDataFrame_k10,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
        legend.position = "none") +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=10")

## -- k = 14 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k14 <- terms(topicModel_k14, 5)
topicNames_k14 <- apply(top5termsPerTopic_k14, 2, paste, collapse = " ")

# get topic proportions from example documents
topicProportions_k14 <- theta_k14[docIds,]
colnames(topicProportions_k14) <- topicNames_k14

vizDataFrame_k14 <- melt(cbind(data.frame(topicProportions_k14),
                                document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k14 <- ggplot(data = vizDataFrame_k14,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
```



```

    legend.position = "none") +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=14")

## -- k = 16 -- ##
# name each topic based on 1st 5 terms in topic
top5termsPerTopic_k16 <- terms(topicModel_k16, 5)
topicNames_k16 <- apply(top5termsPerTopic_k16, 2, paste, collapse = " ")

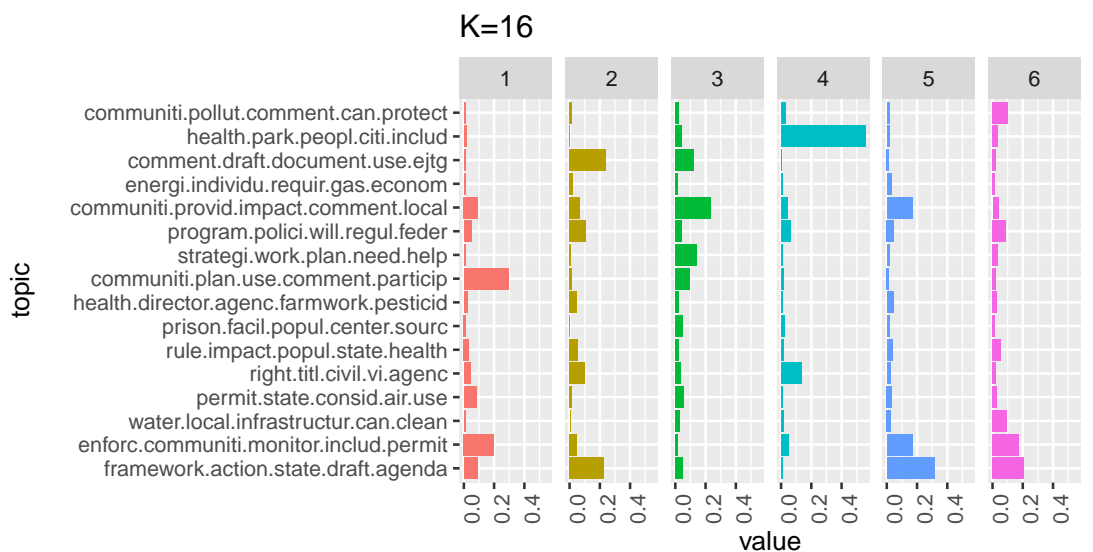
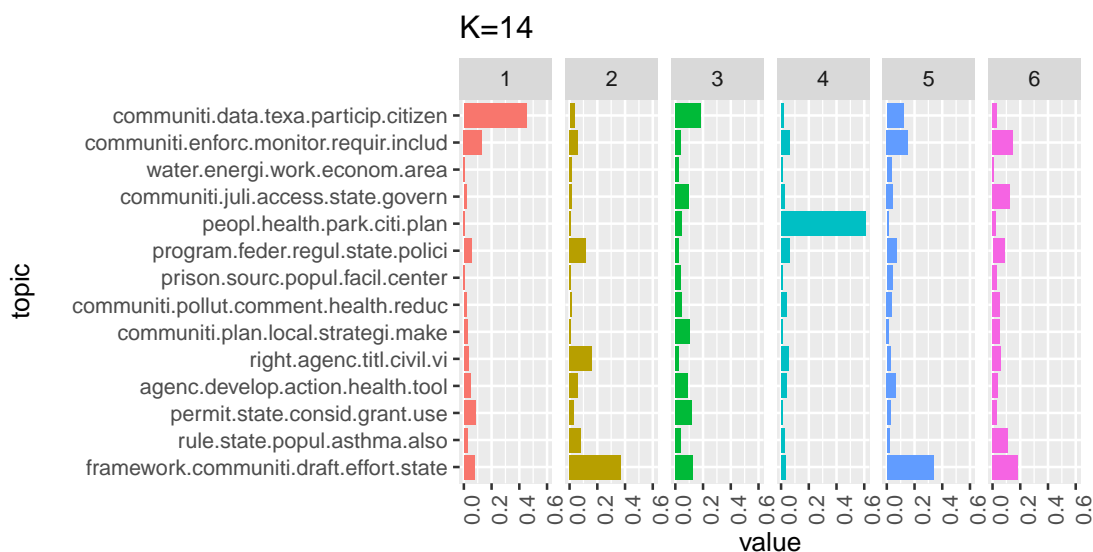
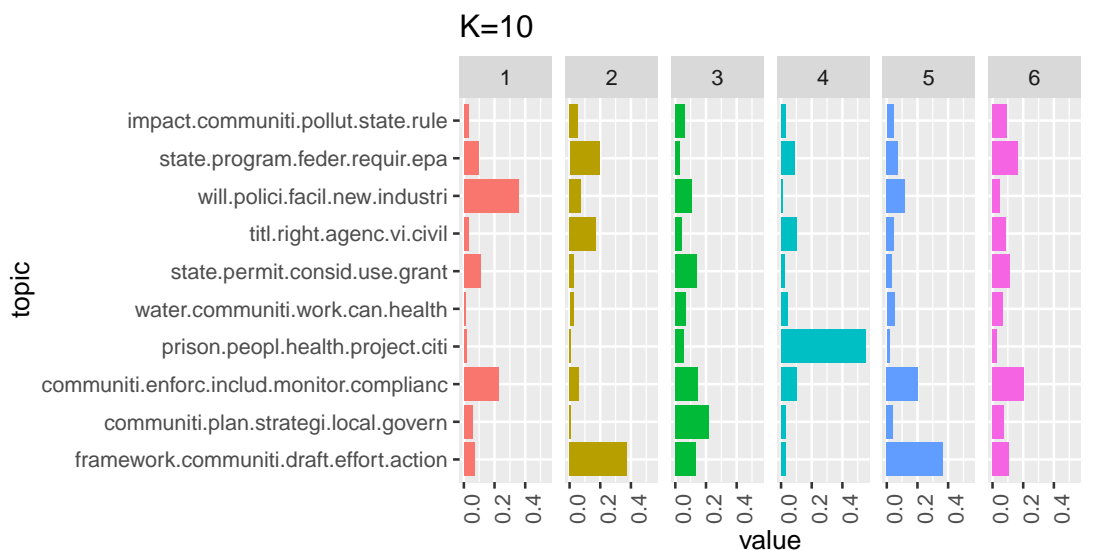
# get topic proportions from example documents
topicProportions_k16 <- theta_k16[docIds,]
colnames(topicProportions_k16) <- topicNames_k16

vizDataFrame_k16 <- melt(cbind(data.frame(topicProportions_k16),
                                   document = factor(1:N)),
                        variable.name = "topic",
                        id.vars = "document")

plot_k16 <- ggplot(data = vizDataFrame_k16,
                  aes(x = topic,
                      y = value,
                      fill = document),
                  ylab = "proportion") +
  geom_bar(stat="identity") +
  theme(axis.text.x = element_text(angle = 90, hjust = 1),
        legend.position = "none") +
  coord_flip() +
  facet_wrap(~ document, ncol = N) +
  labs(title = "K=16")

```

```
plot_k10 / plot_k14 / plot_k16
```



```

svd_tsne <- function(x){
  tsne(svd(x)$u)
}

json_k10 <- createJSON(
  phi = tmResult_k10$terms,
  theta = tmResult_k10$topics,
  doc.length = rowSums(doc_feature_matrix),
  vocab = colnames(doc_feature_matrix),
  term.frequency = colSums(doc_feature_matrix),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))

json_k14 <- createJSON(
  phi = tmResult_k14$terms,
  theta = tmResult_k14$topics,
  doc.length = rowSums(doc_feature_matrix),
  vocab = colnames(doc_feature_matrix),
  term.frequency = colSums(doc_feature_matrix),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))

json_k16 <- createJSON(
  phi = tmResult_k16$terms,
  theta = tmResult_k16$topics,
  doc.length = rowSums(doc_feature_matrix),
  vocab = colnames(doc_feature_matrix),
  term.frequency = colSums(doc_feature_matrix),
  mds.method = svd_tsne,
  plot.opts = list(xlab = "", ylab = ""))

serVis(json_k16)

```

Picking k values to check Based on the `FindTopicsNumber()` optimization metrics, 9, 10, 14, 16, and 20 are good candidates for k (the number of topics). Based on the plots, 10 is the highest maximum and around where the minimization plot begins to decrease more slowly. In the maximization plot, 14 is the next highest maximum peak and 16 is where the minimization plot begins to level out. We know there are 9 priority areas for the EPA and as discussed in the lab, there are 7 additional topics identified in the EPA's response to the public comments. So this would be a theoretical argument for 16 being the best k value. Based on all of the above, I chose to run three models with k equal to 10, 14, and 16.

Picking the best k value Based on the LDAvis and looking at the graphs of the top terms in each topic I would move forward with using 14 as the overall best value for k. When looking at the top term plots for k=16 as compared to k=14, there appeared to begin to be topics that were related to each other such as topics 14 and 15 in the plot that both reference community actions and participation with government. In looking at the LDAvis, the size of the circles for each topic appear to be evenly sized for k=14. Whereas, for k=16 there are more topics with smaller circles (aka less topic distribution).